

## Imitation Learning

As discussed in the previous chapter, the goal of reinforcement learning is to determine closed-loop control policies that result in the maximization of an accumulated reward, and RL algorithms are generally classified as either model-based or model-free. In both cases it is generally assumed that the reward function is known, and both typically rely on collecting system data to either update a learned model (model-based), or directly update a learned value function or policy (model-free).

While successful in many settings, these approaches to RL also suffer from several drawbacks. First, determining an appropriate reward function that can accurately represent the true performance objectives can be challenging<sup>1</sup>. Second, rewards may be *sparse*, which makes the learning process expensive in terms of both the required amount of data and in the number of failures that may be experienced when exploring with a suboptimal policy<sup>2</sup>. This chapter introduces the *imitation learning* approach to RL, where a reward function is not assumed to be known *a priori* but rather it is assumed the reward function is described implicitly through expert demonstrations.

### Imitation Learning

The formulation of the imitation learning problem is quite similar to the RL problem formulation from the previous chapter. The main difference is that instead of leveraging an explicit reward function  $r_t = R(\mathbf{x}_t, \mathbf{u}_t)$  it will be assumed that a set of demonstrations from an expert are provided.

#### 10.1 Problem Formulation

It will be assumed that the system is a Markov Decision Process (MDP) with a state  $\mathbf{x}$  and control input  $\mathbf{u}$ , and the set of admissible states and controls are denoted as  $\mathcal{X}$  and  $\mathcal{U}$ . The system dynamics are expressed by the probabilistic transition model:

$$p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{u}_{t-1}), \quad (10.1)$$

<sup>1</sup> RL agents can sometimes learn how to exploit a reward function without actually producing the desired behavior. This is commonly referred to as *reward hacking*. Consider training an RL agent with a reward for each piece of trash collected. Rather than searching the area to find more trash (the desired behavior), the agent may decide to throw the trash back onto the ground and pick it up again!

<sup>2</sup> This issue of sparse rewards is less relevant if data is cheap, for example when training in simulation.

The field of RL often uses  $\mathbf{s}$  to express the state and  $\mathbf{a}$  to represent an action, but  $\mathbf{x}$  and  $\mathbf{u}$  will be used here for consistency with previous chapters.

which is the conditional probability distribution over  $x_t$ , given the previous state and control. As in the previous chapter, the goal is to define a *policy*  $\pi$  that defines the closed-loop control law<sup>3</sup>:

$$u_t = \pi(x_t). \quad (10.2)$$

The primary difference in formulation from the previous RL problem is that we do not have access to the reward function, and instead we have access to a set of expert demonstrations where each demonstration  $\xi$  consists of a sequence of state-control pairs:

$$\xi = \{(x_0, u_0), (x_1, u_1), \dots\}, \quad (10.3)$$

which are drawn from the expert policy  $\pi^*$ . The imitation learning problem is therefore to determine a policy  $\pi$  that imitates the expert policy  $\pi^*$ :

**Definition 10.1.1** (Imitation Learning Problem). *For a system with transition model (10.1) with states  $x \in \mathcal{X}$  and controls  $u \in \mathcal{U}$ , the imitation learning problem is to leverage a set of demonstrations  $\Xi = \{\xi_1, \dots, \xi_D\}$  from an expert policy  $\pi^*$  to find a policy  $\hat{\pi}^*$  that imitates the expert policy.*

There are generally two approaches to imitation learning: the first is to directly learn how to imitate the expert's policy and the second is to indirectly imitate the policy by instead learning the expert's reward function. This chapter will first introduce two classical approaches to imitation learning (*behavior cloning* and the *DAgger* algorithm) that focus on directly imitating the policy. Then a set of approaches for learning the expert's reward function will be discussed, which is commonly referred to as *inverse reinforcement learning*. The chapter will then conclude with a couple of short discussions into related topics on learning from experts (e.g. through comparisons or physical feedback) as well as on interaction-aware control.

## 10.2 Behavior Cloning

Behavior cloning approaches use a set of expert demonstrations  $\xi \in \Xi$  to determine a policy  $\pi$  that imitates the expert. This can be accomplished through supervised learning techniques, where the difference between the learned policy and expert demonstrations are minimized with respect to some metric. Concretely, the goal is to solve the optimization problem:

$$\hat{\pi}^* = \arg \min_{\pi} \sum_{\xi \in \Xi} \sum_{x \in \xi} L(\pi(x), \pi^*(x)),$$

where  $L$  is the cost function<sup>4</sup>,  $\pi^*(x)$  is the expert's action for at the state  $x$ , and  $\hat{\pi}^*$  is the approximated policy.

However this approach may not yield very good performance since the learning process is only based on a set of samples provided by the expert. In many cases these expert demonstrations will not be uniformly sampled across the

<sup>3</sup> This chapter will consider a stationary policy for simplicity.

<sup>4</sup> Different loss functions could include  $p$ -norms (e.g. Euclidean norm) or  $f$ -divergences (e.g. KL divergence) depending on the form of the policy.

entire state space and therefore it is likely that the learned policy will perform poorly when not close to states found in  $\zeta$ . This is particularly true when the expert demonstrations come from a *trajectory* of sequential states and actions, such that the *distribution* of the sampled states  $x$  in the dataset is defined by the expert policy. Then, when an estimated policy  $\hat{\pi}^*$  is used in practice it produces its own distribution of states that will be visited, which will likely not be the same as in the expert demonstrations! This distributional mismatch leads to compounding errors, which is a major challenge in imitation learning.

### 10.3 DAgger: Dataset Aggregation

One straightforward idea for addressing the issue of distributional mismatch in states seen under the expert policy and the learned policy is to simply collect new expert data as needed<sup>5</sup>. In other words, when the learned policy  $\hat{\pi}^*$  leads to states that aren't in the expert dataset just query the expert for more information! The behavioral cloning algorithm that leverages this idea is known as DAgger<sup>6</sup> (Dataset Aggregation).

<sup>5</sup> Assuming the expert can be queried on demand.

<sup>6</sup> S. Ross, G. Gordon, and D. Bagnell. "A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning". In: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*. 2011, pp. 627–635

---

#### Algorithm 1: DAgger: Dataset Aggregation

---

**Data:**  $\pi^*$

**Result:**  $\hat{\pi}^*$

$\mathcal{D} \leftarrow \emptyset$

Initialize  $\hat{\pi}$

**for**  $i = 1$  **to**  $N$  **do**

$\pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}$

    Rollout policy  $\pi_i$  to sample trajectory  $\tau = \{x_0, x_1, \dots\}$

    Query expert to generate dataset  $\mathcal{D}_i = \{(x_0, \pi^*(x_0)), (x_1, \pi^*(x_1)), \dots\}$

    Aggregate datasets,  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$

    Retrain policy  $\hat{\pi}$  using aggregated dataset  $\mathcal{D}$

**return**  $\hat{\pi}$

---

As can be seen in Algorithm 1, this approach iteratively improves the learned policy by collecting additional data from the expert. This is accomplished by rolling out the current learned policy for some number of time steps and then asking the expert what actions they would have taken at each step along that trajectory. Over time this process drives the learned policy to better approximate the true policy and reduce the incidence of distributional mismatch. One disadvantage to the approach is that at each step the policy needs to be retrained, which may be computationally inefficient.

### 10.4 Inverse Reinforcement Learning

Approaches that learn policies to imitate expert actions can be limited by several factors:

1. Behavior cloning provides no way to understand the underlying reasons for the expert behavior (no reasoning about outcomes or intentions).
2. The “expert” may actually be suboptimal<sup>7</sup>.
3. A policy that is optimal for the expert may not be optimal for the agent if they have different dynamics, morphologies, or capabilities.

An alternative approach to behavioral cloning is to reason about and try to learn a representation of the underlying reward function  $R$  that the expert was using to generate its actions. By learning the expert’s intent, the agent can potentially outperform the expert or adjust for differences in capabilities<sup>8</sup>. This approach (learning reward functions) is known as *inverse reinforcement learning*.

Inverse RL approaches assume a specific parameterization of the reward function, and in this section the fundamental concepts will be presented by parameterizing the reward as a linear combination of (nonlinear) features:

$$R(\mathbf{x}, \mathbf{u}) = \mathbf{w}^T \phi(\mathbf{x}, \mathbf{u}),$$

where  $\mathbf{w} \in \mathbb{R}^n$  is a weight vector and  $\phi(\mathbf{x}, \mathbf{u}) : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}^n$  is a feature map. For a given feature map  $\phi$ , the goal of inverse RL can be simplified to determining the weights  $\mathbf{w}$ . Recall from the previous chapter on RL that the total (discounted) reward under a policy  $\pi$  is defined for a time horizon  $T$  as:

$$V_T^\pi(\mathbf{x}) = E \left[ \sum_{t=0}^{T-1} \gamma^t R(\mathbf{x}_t, \pi(\mathbf{x}_t)) \mid \mathbf{x}_0 = \mathbf{x} \right].$$

Using the reward function  $R(\mathbf{x}, \mathbf{u}) = \mathbf{w}^T \phi(\mathbf{x}, \mathbf{u})$  this value function can be expressed as:

$$V_T^\pi(\mathbf{x}) = \mathbf{w}^T \mu(\pi, \mathbf{x}), \quad \mu(\pi, \mathbf{x}) = E_\pi \left[ \sum_{t=0}^{T-1} \gamma^t \phi(\mathbf{x}_t, \pi(\mathbf{x}_t)) \mid \mathbf{x}_0 = \mathbf{x} \right],$$

where  $\mu(\pi, \mathbf{x})$  is defined by an expectation over the trajectories of the system under policy  $\pi$  (starting from state  $\mathbf{x}$ ) and is referred to as the *feature expectation*<sup>9</sup>. One insight that can now be leveraged is that by definition the optimal expert policy  $\pi^*$  will always produce a greater value function:

$$V_T^{\pi^*}(\mathbf{x}) \geq V_T^\pi(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, \quad \forall \pi,$$

which can be expressed in terms of the feature expectation as:

$$\mathbf{w}^{*T} \mu(\pi^*, \mathbf{x}) \geq \mathbf{w}^{*T} \mu(\pi, \mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, \quad \forall \pi. \quad (10.4)$$

Theoretically, identifying the vector  $\mathbf{w}^*$  associated with the expert policy can be accomplished by finding a vector  $\mathbf{w}$  that satisfies this condition. However this can potentially lead to ambiguities! For example, the choice  $\mathbf{w} = 0$  satisfies this condition trivially! In fact, reward ambiguity is one of the main challenges associated with inverse reinforcement learning<sup>10</sup>. The algorithms discussed in the following chapters will propose techniques for alleviating this issue.

<sup>7</sup> Although the discussion of inverse RL in this section will also assume the expert is optimal, there exist approaches to remove this assumption.

<sup>8</sup> Learned reward representations can potentially generalize across different robot platforms that tackle similar problems!

<sup>9</sup> Feature expectations are often computed using a Monte Carlo technique (e.g. using the set of demonstrations for the expert policy).

<sup>10</sup> A. Ng and S. Russell. “Algorithms for Inverse Reinforcement Learning”. In: *Proceedings of the Seventeenth International Conference on Machine Learning*. 2000, pp. 663–670

### 10.4.1 Apprenticeship Learning

The apprenticeship learning<sup>11</sup> algorithm attempts to avoid some of the problems with reward ambiguity by leveraging an additional insight from condition (10.4). Specifically, the insight is that it doesn't matter how well  $w^*$  is estimated as long as a policy  $\pi$  can be found that *matches the feature expectations*. Mathematically, this conclusion is derived by noting that:

$$\|\mu(\pi, x) - \mu(\pi^*, x)\|_2 \leq \epsilon \implies |w^T \mu(\pi, x) - w^T \mu(\pi^*, x)| \leq \epsilon$$

for any  $w$  as long as  $\|w\|_2 \leq 1$ . In other words, as long as the feature expectations can be matched then the performance will be as good as the expert *even if the vector  $w$  does not match  $w^*$* . Another practical aspect to the approach is that it will be assumed that the initial state  $x_0$  is drawn from a distribution  $D$  such that the value function is also considered in expectation as:

$$E_{x_0 \sim D}[V_T^\pi(x_0)] = w^T \mu(\pi), \quad \mu(\pi) = E_\pi \left[ \sum_{t=0}^{T-1} \gamma^t \phi(x_t, \pi(x_t)) \right].$$

This is useful to avoid having to consider all  $x \in \mathcal{X}$  when matching features<sup>12</sup>.

To summarize, the goal of the apprenticeship learning approach is to find a policy  $\pi$  that matches the feature expectations with respect to the expert policy (i.e. makes  $\mu(\pi)$  as similar as possible to  $\mu(\pi^*)$ )<sup>13</sup>. This is accomplished through Algorithm 2, which uses an iterative approach to finding better policies.

<sup>11</sup> P. Abbeel and A. Ng. "Apprenticeship Learning via Inverse Reinforcement Learning". In: *Proceedings of the Twenty-First International Conference on Machine Learning*. 2004

<sup>12</sup> Trying to find a policy that matches features for every possible starting state  $x$  is likely intractable or even infeasible.

<sup>13</sup> See Example 10.4.1 for an example of why matching features is intuitively useful.

---

#### Algorithm 2: Apprenticeship Learning

---

**Data:**  $\mu(\pi^*), \epsilon$

**Result:**  $\hat{\pi}^*$

Initialize policy  $\pi_0$

**for**  $i = 1$  **to**  $\dots$  **do**

    Compute  $\mu(\pi_{i-1})$  (or approximate via Monte Carlo)

    Solve problem (10.5) with policies  $\{\pi_0, \dots, \pi_{i-1}\}$  to compute  $w_i$  and  $t_i$

$$(w_i, t_i) = \arg \max_{w, t},$$

$$\text{s.t. } w^T \mu(\pi^*) \geq w^T \mu(\pi) + t, \quad \forall \pi \in \{\pi_0, \dots, \pi_{i-1}\},$$

$$\|w\|_2 \leq 1.$$

(10.5)

**if**  $t_i \leq \epsilon$  **then**

$\hat{\pi}^* \leftarrow$  best feature matching policy from  $\{\pi_0, \dots, \pi_{i-1}\}$

**return**  $\hat{\pi}^*$

    Use RL to find an optimal policy  $\pi_i$  for reward function defined by  $w_i$

---

To better understand this algorithm it is useful to further examine the optimization problem (10.5)<sup>14</sup>. Suppose that instead of making  $w$  a decision variable

<sup>14</sup> This problem can be thought of as an inverse RL problem that is seeking to find the reward function vector  $w$  such that the expert *maximally outperforms* the other policies.

it was actually fixed, then the resulting optimization would be:

$$\begin{aligned} t^*(\mathbf{w}) &= \max_t t, \\ \text{s.t. } \mathbf{w}^T \mu(\pi^*) &\geq \mathbf{w}^T \mu(\pi) + t, \quad \forall \pi \in \{\pi_0, \pi_1, \dots\}, \end{aligned}$$

which is essentially computing the smallest performance loss among the candidate policies  $\{\pi_0, \pi_1, \dots\}$  with respect to the expert policy, *assuming the reward function weights are  $\mathbf{w}$* . If  $\mathbf{w}$  was known, then if  $t^*(\mathbf{w}) \leq \epsilon$  it would guarantee that one of the candidate policies would effectively perform as well as the expert.

Since  $\mathbf{w}$  is not known, the actual optimization problem (10.5) maximizes the smallest performance loss across *all vectors  $\mathbf{w}$  with  $\|\mathbf{w}\|_2 \leq 1$* . Therefore, if  $t_i \leq \epsilon$  (i.e. the termination condition in Algorithm 2), then there must be a candidate policy whose performance loss is small *for all possible choices of  $\mathbf{w}$* ! In other words, there is a candidate policy that matches feature expectations well enough that good performance can be guaranteed without assuming the reward function is known, and without attempting to estimate the reward accurately.

**Example 10.4.1** (Apprenticeship Learning vs. Behavioral Cloning). Consider a problem where the goal is to drive a car across a city in as short of time as possible. In the imitation learning formulation it is assumed that the reward function is not known, but that there is an expert who shows how to drive across the city (i.e. what routes to take). A behavioral cloning approach would simply try to mimic the actions taken by the expert, such as memorizing that whenever the agent is at a particular intersection it should turn right. Of course this approach is not robust when at intersections that the expert never visited!

The apprenticeship learning approach tries to avoid the inefficiency of behavioral cloning by instead identifying features of the expert's trajectories that are more generalizable, and developing a policy that experiences the same feature expectations as the expert. For example it could be more efficient to notice that the expert takes routes without stop signs, or routes with higher speed limits, and then try to find policies that also seek out those features!

#### 10.4.2 Maximum Margin Planning

The maximum margin planning approach<sup>15</sup> uses an optimization-based approach to computing the reward function weights  $\mathbf{w}$  that is very similar to (10.5) but with some additional flexibility. In its most standard form the MMP optimization is:

$$\begin{aligned} \hat{\mathbf{w}}^* &= \arg \min_{\mathbf{w}} \|\mathbf{w}\|_2^2, \\ \text{s.t. } \mathbf{w}^T \mu(\pi^*) &\geq \mathbf{w}^T \mu(\pi) + 1, \quad \forall \pi \in \{\pi_0, \pi_1, \dots\}. \end{aligned}$$

Again this problem computes the reward function vector  $\mathbf{w}$  such that the expert policy *maximally outperforms* the policies in the set  $\{\pi_0, \pi_1, \dots\}$ .

<sup>15</sup> N. Ratliff, J. A. Bagnell, and M. Zinkevich. "Maximum Margin Planning". In: *Proceedings of the 23rd International Conference on Machine Learning*. 2006, pp. 729–736

However the formulation is also improved in two ways: it adds a slack term to account for potential expert suboptimality and it adds a similarity function that gives more “margin” to policies that are dissimilar to the expert policy. This new formulation is:

$$\begin{aligned} \hat{w}^* &= \arg \min_{w,v} \|w\|_2^2 + Cv, \\ \text{s.t. } w^T \mu(\pi^*) &\geq w^T \mu(\pi) + m(\pi^*, \pi) - v, \quad \forall \pi \in \{\pi_0, \pi_1, \dots\}, \end{aligned} \quad (10.6)$$

where  $v$  is a slack variable that can account for expert suboptimality,  $C > 0$  is a hyperparameter that is used to penalize the amount of assumed suboptimality, and  $m(\pi^*, \pi)$  is a function that quantifies how dissimilar two policies are.

One example of where this formulation is advantageous over the apprenticeship learning formulation (10.5) is when the expert is suboptimal. In this case it is possible that there is no  $w$  that makes the expert policy outperform all other policies, such that the optimization (10.5) returns  $w_i = 0$  and  $t_i = 0$  (which is obviously not the appropriate solution). Alternatively the slack variables in the MMP formulation allow for a reasonable  $w$  to be computed.

### 10.4.3 Maximum Entropy Inverse Reinforcement Learning

While the apprenticeship learning approach shows that matching feature counts is a necessary and sufficient condition to ensure a policy performs as well as an expert, it also has some ambiguity (similar to the reward weight ambiguity problem discussed before). This ambiguity is associated with the fact that there could be different policies that lead to the same feature expectations!

This issue can also be thought of in a slightly more intuitive way in terms of distributions over trajectories. Specifically, a policy  $\pi$  induces a distribution over trajectories<sup>16</sup>  $\tau = \{(x_0, \pi(x_0)), (x_1, \pi(x_1)), \dots\}$  that is denoted as  $p_\pi(\tau)$ . The feature expectations can be rewritten in terms of this distribution as:

$$\mu(\pi) = E_\pi[f(\tau)] = \int p_\pi(\tau) f(\tau) d\tau,$$

where  $f(\tau) = \sum_{t=0}^{T-1} \gamma^t \phi(x_t, \pi(x_t))$ . Now suppose a policy  $\pi$  was found that matched feature expectations<sup>17</sup> with an expert policy  $\pi^*$  such that:

$$\int p_\pi(\tau) f(\tau) d\tau = \int p_{\pi^*}(\tau) f(\tau) d\tau.$$

Crucially this condition is not sufficient to guarantee that  $p_\pi(\tau) = p_{\pi^*}(\tau)$  (which would be ideal). In fact, the distribution  $p_\pi(\tau)$  could also have an *arbitrary preference* for some paths that is *unrelated* to the feature matching objective.

The main idea in the maximum entropy inverse RL approach<sup>18</sup> is to not only match the feature expectations, but also remove ambiguity in the path distribution  $p_\pi(\tau)$  by trying to make  $p_\pi(\tau)$  as *broadly uncommitted as possible*. In other words, find a policy that matches feature expectations but otherwise has no additional path preferences. This concept is known as the maximum entropy principle<sup>19</sup>.

<sup>16</sup> This distribution can be visualized as a set of paths generated by simulating the system many times with policy  $\pi$  (i.e. using a Monte Carlo method).

<sup>17</sup> For example by using apprenticeship learning.

<sup>18</sup> B. D. Ziebart et al. “Maximum Entropy Inverse Reinforcement Learning”. In: *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence*. 2008, pp. 1433–1438

<sup>19</sup> A maximum entropy distribution can be thought of as the *least informative distribution* of a class of distribution. This is useful in situations where it is undesirable to encode unintended prior information.

The maximum entropy IRL approach finds a minimally preferential, feature expectation matching distribution by solving the optimization problem:

$$\begin{aligned}
p^*(\tau) &= \arg \max_p \int -p(\tau) \log p(\tau) d\tau, \\
\text{s.t. } &\int p(\tau) f(\tau) d\tau = \int p_{\pi^*}(\tau) f(\tau) d\tau, \\
&\int p(\tau) d\tau = 1, \\
&p(\tau) \geq 0, \quad \forall \tau,
\end{aligned} \tag{10.7}$$

where the objective is the mathematical definition of a distribution's entropy, the first constraint requires feature expectation matching, and the remaining constraints ensure that  $p(\tau)$  is a valid probability distribution. It turns out that the solution to this problem has the exponential form:

$$p^*(\tau, \lambda) = \frac{1}{Z(\lambda)} e^{\lambda^T f(\tau)}, \quad Z(\lambda) = \int e^{\lambda^T f(\tau)} d\tau,$$

where  $Z(\lambda)$  normalizes the distribution, and where  $\lambda$  must be chosen such that the feature expectations match:

$$\int p^*(\tau, \lambda) f(\tau) d\tau = \int p_{\pi^*}(\tau) f(\tau) d\tau.$$

In other words the maximum entropy IRL approach tries to find a distribution parameterized by  $\lambda$  that match features, but also requires that the distribution  $p^*(\tau, \lambda)$  belong to the exponential family.

To determine the value of  $\lambda$  that matches features, it is assumed that the expert also selects trajectories with high reward with exponentially higher probability:

$$p_{\pi^*}(\tau) \propto e^{w^* T f(\tau)},$$

and therefore ideally  $\lambda = w^*$ . Of course  $w^*$  (and more generally  $p_{\pi^*}(\tau)$ ) are not known, and therefore a maximum likelihood estimation approach is used to compute  $\lambda$  to best approximate  $w^*$  based on the sampled expert demonstrations<sup>20</sup>.

In particular, an estimate  $\hat{w}^*$  of the reward weights is computed from the expert demonstrations  $\Xi = \{\xi_0, \xi_1, \dots\}$  (which each demonstration  $\xi_i$  is a trajectory) by solving the maximum likelihood problem:

$$\begin{aligned}
\hat{w}^* &= \arg \max_{\lambda} \prod_{\xi_i \in \Xi} p^*(\xi_i, \lambda), \\
&= \arg \max_{\lambda} \sum_{\xi_i \in \Xi} \lambda^T f(\xi_i) - \log Z(\lambda),
\end{aligned}$$

which can be solved using a gradient descent algorithm where the gradient is computed by:

$$\nabla_{\lambda} J(\lambda) = \sum_{\xi_i \in \Xi} f(\xi_i) - E_{\tau \sim p^*(\tau, \lambda)}[f(\tau)].$$

<sup>20</sup> By assuming the expert policy is also exponential, the maximum likelihood estimate is theoretically *consistent* (i.e.  $\lambda \rightarrow w^*$  as the number of demonstrations approaches infinity).



The first term of this gradient is easily computable since the expert demonstrations are known, and the second term can be approximated through Monte Carlo sampling. However, this Monte Carlo sampling estimate is based on sampling trajectories from the distribution  $p^*(\tau, \lambda)$ . This leads to an iterative algorithm:

1. Initialize  $\lambda$  and collect the set of expert demonstrations  $\Xi = \{\xi_0, \xi_1, \dots\}$ .
2. Compute the optimal policy<sup>21</sup>  $\pi_\lambda$  with respect to the reward function with  $w = \lambda$ .
3. Using the policy  $\pi_\lambda$ , sample trajectories of the system and compute an approximation of  $E_{\tau \sim p^*(\tau, \lambda)}[f(\tau)]$ .
4. Perform a gradient step on  $\lambda$  to improve the maximum likelihood cost.
5. Repeat until convergence.

<sup>21</sup> For example through traditional RL methods.

To summarize, the maximum entropy inverse reinforcement learning approach identifies a distribution over trajectories that matches feature expectations with the expert, but by restricting the distribution to belong to the exponential family ensures that spurious preferences (path preferences not motivated by feature matching) are not introduced. Additionally, this distribution over trajectories is parameterized by a value that is an estimate of the reward function weights.

## 10.5 Learning From Comparisons and Physical Feedback

Both behavioral cloning and inverse reinforcement learning approaches rely on expert demonstrations of behavior. However in some practical scenarios it may actually be difficult for the expert to provide complete/quality demonstrations. For example it has been shown<sup>22</sup> that when humans are asked to demonstrate good driving behavior in simulation they retroactively think their behavior was too aggressive! As another example, if a robot has a high-dimensional control or state space it could be difficult for the expert to specify the full high-dimensional behavior. Therefore another interesting question in imitation learning is to find a way to learn from alternative data sources besides complete demonstrations.

<sup>22</sup> C. Basu et al. "Do You Want Your Autonomous Car to Drive Like You?" In: *12th ACM/IEEE International Conference on Human-Robot Interaction*. 2017, PP. 417-425

### 10.5.1 Learning from Comparisons

One alternative approach is to use *pairwise comparisons*<sup>23</sup>, where an expert is shown two different behaviors and then asked to rank which behavior is better. Through repeated queries it is possible to converge to an understanding of the underlying reward function. For example, suppose two trajectories  $\tau_A$  and  $\tau_B$  are shown to an expert and that trajectory  $\tau_A$  is preferred. Then assuming that the reward function is:

$$R(\tau) = w^T f(\tau),$$

<sup>23</sup> D. Sadigh et al. "Active Preference-Based Learning of Reward Functions". In: *Robotics: Science and System*. 2017

where  $f(\tau)$  are the collective feature counts (same as in Section 10.4), this comparison can be used to conclude that:

$$\mathbf{w}^T f(\tau_A) > \mathbf{w}^T f(\tau_B).$$

In other words, this comparison has split the space of possible reward weights  $\mathbf{w}$  in half through the hyperplane:

$$(f(\tau_A) - f(\tau_B))^T \mathbf{w} = 0.$$

By continuously querying the expert with new comparisons<sup>24</sup>, the space of possible reward weights  $\mathbf{w}$  will continue to shrink until a good estimate of  $\mathbf{w}^*$  can be made. In practice the expert decision may be a little noisy and therefore the hyperplanes don't define hard cutoffs, but rather can be used to "weight" the possible reward vectors  $\mathbf{w}$ .

<sup>24</sup> The types of comparisons shown can be selectively chosen to maximally split the remaining space of potential  $\mathbf{w}$  in order to minimize the total number of expert queries that are required.

### 10.5.2 Learning from Physical Feedback

Another alternative to learning from complete expert demonstrations is to simply allow the expert to physically interact with the robot to correct for undesirable behavior<sup>25</sup>. In this approach, a physical interaction (i.e. a correction) is assumed to occur when the robot takes actions that result in a lower reward than the expert's action.

For a reward function of the form  $R(\mathbf{x}, \mathbf{u}) = \mathbf{w}^T \phi(\mathbf{x}, \mathbf{u})$  the robot maintains an estimate of the reward weights  $\hat{\mathbf{w}}^*$  and the expert is assumed to have act according to a true set of optimal weights  $\mathbf{w}^*$ . Suppose the robot's policy, which is based on the estimated reward function with weights  $\hat{\mathbf{w}}^*$ , yields a trajectory  $\tau_R$ . Then, if the expert physically interacts with the robot to make a correction the resulting actual trajectory  $\tau_H$  is assumed to satisfy:

$$\mathbf{w}^{*T} f(\tau_H) \geq \mathbf{w}^{*T} f(\tau_R),$$

which simply states that the reward of the new trajectory is higher. This insight is then leveraged in a maximum a posteriori approach for updating the estimate  $\hat{\mathbf{w}}^*$  after each interaction. Specifically, this update takes the form:

$$\hat{\mathbf{w}}^* \leftarrow \hat{\mathbf{w}}^* + \beta(f(\tau_H) - f(\tau_R)),$$

where  $\beta > 0$  is a scalar step size. The robot then uses the new estimate to change its policy, and the process iterates. Note that this idea yields an approach that is similar to the concept of matching feature expectations from inverse reinforcement learning, except that the approach is iterative rather than requiring a batch of complete expert demonstrations.

### 10.6 Interaction-aware Control and Intent Inference

Yet another interesting problem in robot autonomy arises when robots and humans are interacting to accomplish shared or individual goals. Many classical

<sup>25</sup> A. Bajcsy et al. "Learning Robot Objectives from Physical Human Interaction". In: *Proceedings of the 1st Annual Conference on Robot Learning*. 2017, pp. 217–226

examples of this problem arise in autonomous driving settings, when human-driven vehicles interact with autonomous vehicles in settings such as highway merging or at intersections. While the imitation learning problems from the previous sections are focused on understanding the expert's behavior for the purpose of *imitating* the behavior, in this setting the human's behavior needs to be understood in order to ensure safe interactions. However there is an additional component to understanding interactions: *the robot's behavior can influence the human's behavior*<sup>26</sup>.

### 10.6.1 Interaction-aware Control with Known Human Model

One common approach is to model the interaction between humans and robots as a dynamical system that has a combined state  $\mathbf{x}$ , where the robot controls are denoted  $\mathbf{u}_R$  and the human decisions or inputs are denoted as  $\mathbf{u}_H$ . The transition model is therefore defined as:

$$p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{u}_{R,t-1}, \mathbf{u}_{H,t-1}).$$

In other words the interaction dynamics evolve according to the actions taken by both the robot and the human. In this interaction the robot's reward function is denoted as  $R_R(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H)$  and the human's reward function is denoted as  $R_H(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H)$ , which are both functions of the combined state and both agent's actions<sup>27</sup>.

Under the assumption that both the robot and the human act optimally<sup>28</sup> with respect to their cost functions:

$$\begin{aligned} \mathbf{u}_R^*(\mathbf{x}) &= \arg \max_{\mathbf{u}_R} R_R(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H^*(\mathbf{x})), \\ \mathbf{u}_H^*(\mathbf{x}) &= \arg \max_{\mathbf{u}_H} R_H(\mathbf{x}, \mathbf{u}_R^*(\mathbf{x}), \mathbf{u}_H). \end{aligned}$$

Additionally, assuming both reward functions  $R_R$  and  $R_H$  are known<sup>29</sup>, computing  $\mathbf{u}_R^*$  is still extremely challenging due to the two-player game dynamics of the decision making process. However this problem can be made more tractable by modeling it as a *Stackelberg game*, which restricts the two-player game dynamics to a leader-follower structure. Under this assumption it is assumed that the robot is the "leader" and that as the follower the human acts according to:

$$\mathbf{u}_H^*(\mathbf{x}, \mathbf{u}_R) = \arg \max_{\mathbf{u}_H} R_H(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H). \quad (10.8)$$

In other words the human is assumed to see the action taken by the robot *before* deciding on their own action. The robot policy can therefore be computed by solving:

$$\mathbf{u}_R^*(\mathbf{x}) = \arg \max_{\mathbf{u}_R} R_R(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H^*(\mathbf{x}, \mathbf{u}_R)), \quad (10.9)$$

which can be solved using a gradient descent approach. For the gradient descent approach the gradient of:

$$J(\mathbf{x}, \mathbf{u}_R) = R_R(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H^*(\mathbf{x}, \mathbf{u}_R)),$$

<sup>26</sup> It is particularly important in interaction-aware robot control to understand the effects of the robot's actions on the human's behavior. Otherwise the human's could simply be modeled as dynamic obstacles!

<sup>27</sup> While  $R_R$  and  $R_H$  do not have to be the same, choosing  $R_R = R_H$  may be desirable for the robot to achieve human-like behavior.

<sup>28</sup> While not necessarily true, this assumption is important to make the resulting problem formulation tractable to solve in practice.

<sup>29</sup> The reward function  $R_H$  could be approximated using inverse reinforcement learning techniques.

can be computed using the chain rule as:

$$\frac{\partial J}{\partial \mathbf{u}_R} = \frac{\partial R_R}{\partial \mathbf{u}_R} + \frac{\partial R_R}{\partial \mathbf{u}_H^*} \frac{\partial \mathbf{u}_H^*}{\partial \mathbf{u}_R}.$$

Since the reward function  $R_R$  is known the terms  $\partial R_R / \partial \mathbf{u}_R$  and  $\partial R_R / \partial \mathbf{u}_H^*$  can be easily determined. In order to compute the term  $\partial \mathbf{u}_H^* / \partial \mathbf{u}_R$ , which represents how much the robot's actions impact the human's actions, an additional step is required. First, assuming the human acts optimally according to (10.8) the necessary optimality condition is:

$$g(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H^*) = 0, \quad g = \frac{\partial R_H}{\partial \mathbf{u}_H},$$

which for the fixed values of  $\mathbf{x}$  and  $\mathbf{u}_R$  specifies  $\mathbf{u}_H^*$ . Then, by implicitly differentiating this condition with respect to the robot action  $\mathbf{u}_R$ :

$$\frac{\partial g}{\partial \mathbf{u}_R} + \frac{\partial g}{\partial \mathbf{u}_H^*} \frac{\partial \mathbf{u}_H^*}{\partial \mathbf{u}_R} = 0,$$

which can be used to solve for:

$$\frac{\partial \mathbf{u}_H^*}{\partial \mathbf{u}_R}(\mathbf{x}, \mathbf{u}_R, \mathbf{u}_H^*) = - \left( \frac{\partial g}{\partial \mathbf{u}_H^*} \right)^{-1} \frac{\partial g}{\partial \mathbf{u}_R}.$$

Notice that every term in this expression can be computed<sup>30</sup> and therefore it can be substituted into the gradient calculation:

$$\frac{\partial J}{\partial \mathbf{u}_R} = \frac{\partial R_R}{\partial \mathbf{u}_R} - \frac{\partial R_R}{\partial \mathbf{u}_H^*} \left( \frac{\partial g}{\partial \mathbf{u}_H^*} \right)^{-1} \frac{\partial g}{\partial \mathbf{u}_R},$$

which can then be computed as long as it is possible to compute  $\mathbf{u}_H^*(\mathbf{x}, \mathbf{u}_R)$ .

To summarize, one approach to interaction-aware control is to model the interaction as a Stackelberg game, where it is assumed that both the human and the robot act optimally with respect to some reward functions. This formulation of the problem enables the robot to choose actions based on an implicit understanding of how the human will react.

### 10.6.2 Intent Inference

One disadvantage to the approach for interaction-aware control from the previous section is that it assumes the human acts optimally with respect to a *known* reward function. While a reward function could be learned through inverse reinforcement learning, this is not practical for real-world settings where different humans behave differently. Returning to the example of interaction between human drivers and autonomous vehicles, the human could exhibit drastically different behavior depending on whether they have an aggressive or passive driving style. In these settings the problem of *intent inference* focuses on identifying underlying behavioral characteristics that can lead to more accurate behavioral models<sup>31</sup>.

<sup>30</sup> Assuming the human's reward function is known.

<sup>31</sup> This problem can be formulated as a partially observable Markov decision process (POMDP) since the underlying behavioral characteristic is not directly observable, yet influences the system's behavior.

One approach to intent inference<sup>32</sup> is to model the underlying behavioral differences through a set of unknown parameters  $\theta$  which need to be inferred by observing the human's behavior. Mathematically this is expressed by defining the human's reward function  $R_H(x, u_R, u_H, \theta)$  to be a function of  $\theta$ , and assuming the human chooses actions according to:

$$p(u_H | x, u_R, \theta) \propto e^{R_H(x, u_R, u_H, \theta)}.$$

In other words this model assumes the human is exponentially more likely to pick optimal actions<sup>33</sup>, but that they may pick suboptimal actions as well.

The objective of intent inference is therefore to estimate the parameters  $\theta$ , which can be accomplished through Bayesian inference methods. In the Bayesian approach a *probability distribution* over parameters  $\theta$  is updated based on observations. Specifically the belief distribution is denoted as  $b(\theta)$ , and given an observation of the human's actions  $u_H$  the belief distribution is updated as:

$$b_{t+1}(\theta) = \frac{1}{\eta} p(u_{H,t} | x_t, u_{R,t}, \theta) b_t(\theta),$$

where  $\eta$  is a normalizing constant. This Bayesian update is simply taking the prior belief over  $\theta$  and updating the distribution based on the likelihood of observing human action  $u_H$  under that prior. Note that this concept is quite similar to the concepts of inverse reinforcement learning: a set of parameters that describe the human's (experts) behavior are continually updated when new observations of their actions are gathered.

While the robot could sit around and *passively* observe the human act to collect samples for the Bayesian updates, it is often more efficient for the robot to *probe* the human to take interesting actions that are more useful for revealing the intent parameters  $\theta$ . This can be accomplished by choosing the robot's reward function to be:

$$R_R(x, u_R, u_H, \theta) = I(b(\theta), u_R) + \lambda R_{\text{goal}}(x, u_R, u_H, \theta)$$

where  $\lambda > 0$  is a tuning parameter and  $I(b(\theta), u_R)$  denotes a function that quantifies the amount of information gained with respect to the belief distribution from taking action  $u_R$ . In other words the robot's reward is a tradeoff between exploiting the current knowledge of  $\theta$  to accomplish the objective and taking exploratory actions to improve the intent inference. With this robot reward function the robot's actions are chosen to maximize the expected reward:

$$u_R^*(x) = \arg \max_{u_R} E_{\theta}[R_R(x, u_R, u_H, \theta)].$$

To summarize, this robot policy will try to simultaneously accomplish the robot's objective and gather more information to improve the inference of the human's intent (modeled through the parameters  $\theta$ ). In a highway lane changing scenario this type of policy might lead the robot to nudge into the other lane to see if the other car will slow down (passive driving behavior) or try to

<sup>32</sup> D. Sadigh et al. "Planning for cars that coordinate with people: leveraging effects on human actions for planning and active information gathering over human internal state". In: *Autonomous Robots* 42:7 (2018), pp. 1405–1426

<sup>33</sup> This assumption was also used in the Maximum Entropy IRL approach.

block the lane change (aggressive driving behavior). Once the robot has a strong enough belief about the human's behavior it may choose to either complete the lane change or slow down to merge behind the human driver.